

Article

Integrated Analysis of Gene Expression and Protein–Protein Interaction with Tensor Decomposition

Y-H. Taguchi ^{1,*}  and Turki Turki ² ¹ Department of Physics, Chuo University, Tokyo 112-8551, Japan² Department of Computer Science, King Abdulaziz University, Jeddah 21589, Saudi Arabia; tturki@kau.edu.sa

* Correspondence: tag@granular.com; Tel.: +81-3-3817-1791

Abstract: Integration of gene expression (GE) and protein–protein interaction (PPI) is not straightforward because the former is provided as a matrix, whereas the latter is provided as a network. In many cases, genes processed with GE analysis are refined further based on a PPI network or vice versa. This is hardly regarded as a true integration of GE and PPI. To address this problem, we proposed a tensor decomposition (TD)-based method that can integrate GE and PPI prior to any analyses where PPI is also formatted as a matrix to which singular value decomposition (SVD) is applied. Integrated analyses with TD improved the coincidence between vectors attributed to samples and class labels over 27 cancer types retrieved from The Cancer Genome Atlas Program (TCGA) toward five class labels. Enrichment using genes selected with this strategy was also improved with the integration using TD. The PPI network associated with the information on the strength of the PPI can improve the performance than PPI that stores only if the interaction exists in individual pairs. In addition, even restricting genes to the intersection of GE and PPI can improve coincidence and enrichment.

Keywords: tensor decomposition; gene expression; protein–protein interaction; integrated analyses

MSC: 92B05; 68T09



Citation: Taguchi, Y.-H.; Turki, T. Integrated Analysis of Gene Expression and Protein–Protein Interaction with Tensor Decomposition. *Mathematics* **2023**, *11*, 3655. <https://doi.org/10.3390/math11173655>

Academic Editor: Biao Tang

Received: 19 July 2023

Revised: 9 August 2023

Accepted: 21 August 2023

Published: 24 August 2023



Copyright: © 2023 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

The integrated analysis of gene expression (GE) and protein–protein interaction (PPI) has not been studied extensively [1]. This is possibly because of the distinct formats provided for GE and PPI. GE is usually provided in a matrix format, whose columns and rows typically correspond to samples and genes, whereas PPI is usually provided in a network format. Integrating these two formats is not straightforward. Mainly, two approaches were tried: the matrix-based (MB) method and the network-based (NB) method (Figure 1).

Typically, in MB methods, differentially expressed genes (DEGs) are identified and PPI is used to validate or filter DEGs [2–4]. However, these kinds of MB methods are rarely recognized as integrated analyses of GE and PPI since the strategy of further screening DEGs is based on an independent strategy that is common and not restricted to PPI. For example, various enrichment analyses based on previous biological knowledge were often used to screen DEGs. Thus, MB methods are typically unlikely to be regarded as integrated analyses of GE and PPI. In actuality, although [1] reviewed some studies that were regarded as integrated analyses of GE and PPI, most of these studies are not MB but NB methods.

In contrast to MB methods, in NB methods, GE information is mapped onto a network, and modules associated with co-expression genes are selected. These NB approaches integrated with GE are known to improve the performance of simple NB approaches [5–7]. This strategy, where genes embedded in network structures that are further screened based on GE, is more likely to be regarded as integrated analyses of GE and PPI than MB methods because GE is not the only criterion to further screen genes embedded in a network. Many

other criteria, such as the previously mentioned enriched analyses, can be used to further screen genes embedded in a network, and NB methods are not specifically integrated analyses of GE and PPI.

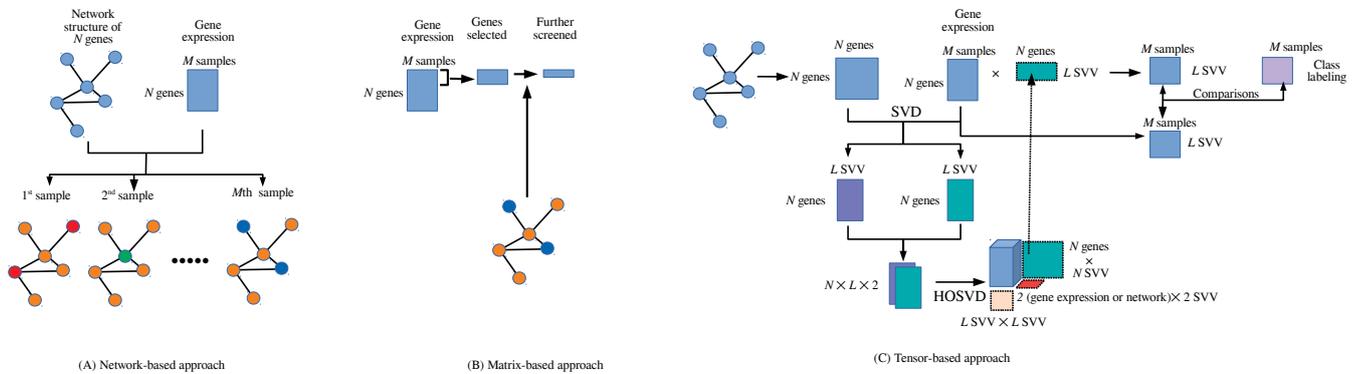


Figure 1. (A) Network-based approach. (B) Matrix-based approach (C). Tensor-based approach (this also shows the analysis flow chart in this study).

The weak point of both the MB and NB methods is obvious. One of two criteria, GE or PPI, inevitably dominates another. In MB approaches, no non-DEGs are selected because only DEGs are passed to be screened by PPI, whereas in NB approaches, no non-network-connected genes are selected because only genes connected within a network are passed to be screened based on whether they are DEGs. Thus, there are inevitable inequalities between GE and PPI in both approaches.

Some studies attempted to avoid the inequality between GE and PPI. For example, ref. [8] tried to equally weight GE and PPI by introducing the Jaccard similarity index (JDC), which is a simple product between the edge clustering coefficient computed from the PPI network and the Jaccard coefficient computed from the GE similarity. Although JDC was superior to other existing methods and removing the inequality between GE and PPI could improve performance, using JDC is restricted to the time course dataset and this type of approach (i.e., trying to equally weight GE and PPI) has only been rarely investigated.

In this paper, we introduce new, tensor-based (TB) approaches to integrate PPI and GE, assuming no priority between GE and PPI. In the TB approach, PPI expressed in a matrix format is once transformed using singular value decomposition into singular value vectors (SVV), which are later bundled with SVVs computed from GE with SVD to generate a tensor to which TD is applied. The TD-generated SVVs attributed to a gene are further used to generate vectors attributed to samples by projecting GE to TD-generated SVVs attributed to a gene. These obtained vectors attributed to samples are tested if they are coincident with class labels (e.g., patients vs. healthy controls) and are selected based on coincidence. Once vectors attributed to samples are selected, then corresponding SVVs attributed to a gene, to which GE is mapped, are used to select DEGs using a previously proposed TD-based unsupervised feature extraction (FE) [9] criterion that was recently improved with optimized standard deviation [10,11] used in Gaussian distribution, which SVVs attributed to genes are supposed to obey in the null hypothesis. Selected genes can be further validated with various enrichment analyses based on previously obtained biological knowledge.

As the scarcity of data hinders the performance of data-driven methods, thereby affecting the biological reliability of selected genes, we proposed a new computational method to integrate GE and PPI and improve the biological significance of obtained results from enrichment analysis as shown from the analysis of cancer data obtained from the GEO database.

2. Materials and Methods

In this paper, to express network structure in a matrix format to be integrated with GE, we employed the simplest network structure representation (NSR) [12], where $n_{ii'} \in \mathbb{R}^{N \times N}$ is 1 only when the node i and i' are connected with each other and otherwise is zero.

The R source code used to perform this analysis is in the supplementary file.

2.1. Integrated Analysis of Matrices and Networks with Tensor

To convert an NSR into a matrix starting with $n_{ii'}$, we applied SVD to $n_{ii'}$ as

$$n_{ii'} = \sum_{\ell} \lambda_{\ell} u_{\ell i} u_{\ell i'} \tag{1}$$

where λ_{ℓ} is a singular value and $u_{\ell i} \in \mathbb{R}^{N \times N}$ are the singular value matrix and the orthogonal matrix.

On the other hand, we apply SVD to $x_{ij} \in \mathbb{R}^{N \times M}$, which represents the gene expression of the i th gene of the j th sample as

$$x_{ij} = \sum_{\ell} \lambda'_{\ell} u'_{\ell i} v_{\ell j} \tag{2}$$

where $u'_{\ell i} \in \mathbb{R}^{M \times N}$ and $v_{\ell j} \in \mathbb{R}^{M \times M}$ (here we assume $M < N$). Then, we generate the tensor $x_{i\ell k} \in \mathbb{R}^{N \times L \times 2}$ using the first $L (< M)$ SVVs as:

$$x_{i\ell k} = \begin{cases} u_{\ell i}, & k = 1 \\ u'_{\ell i}, & k = 2 \end{cases} \tag{3}$$

to which higher order singular value decomposition (HOSVD) is then applied, giving:

$$x_{i\ell k} = \sum_{\ell_1=1}^N \sum_{\ell_2=1}^L \sum_{\ell_3=1}^2 G(\ell_1 \ell_2 \ell_3) \tilde{u}_{\ell_1 i} \tilde{u}_{\ell_2 \ell} \tilde{u}_{\ell_3 k} \tag{4}$$

where $G \in \mathbb{R}^{N \times L \times 2}$ is a core tensor and $\tilde{u}_{\ell_1 i} \in \mathbb{R}^{N \times N}$, $\tilde{u}_{\ell_2 \ell} \in \mathbb{R}^{L \times L}$, and $\tilde{u}_{\ell_3 k} \in \mathbb{R}^{2 \times 2}$ are singular value matrices and orthogonal matrices. Throughout this article, $L = 10$. Next, we project x_{ij} onto $\tilde{u}_{\ell_1 i}$ by

$$\tilde{v}_{\ell_1 j} = \sum_{i=1}^N \tilde{u}_{\ell_1 i} x_{ij} \tag{5}$$

to get vectors attributed to sample, $\tilde{v}_{\ell_1 j}$.

2.2. Comparisons of Coincidence with Class Labels between $v_{\ell j}$ and $\tilde{v}_{\ell_1 j}$

Categorical regression is performed for $v_{\ell j}$ and $\tilde{v}_{\ell_1 j}$ as

$$v_{\ell j} = a_{\ell} + \sum_s b_{\ell s} \delta_{js} \tag{6}$$

$$\tilde{v}_{\ell_1 j} = a'_{\ell_1} + \sum_s b'_{\ell_1 s} \delta_{js} \tag{7}$$

where δ_{js} is 1 only when j belongs to the s th class and otherwise is 0, and $a_{\ell}, b_{\ell s}, a'_{\ell_1}, b'_{\ell_1 s}$ are regression coefficients. p -values are computed by the `lm` function in R. Hereafter, we denote p -values computed using Equations (6) and (7) for the m th dataset as P_{ℓ}^m and \tilde{P}_{ℓ}^m , respectively. The obtained p -values are corrected by the BH criterion [9] using the `p.adjust` function in R with the "BH" option.

If there are M_0 data sets (i.e., $m \leq M_0$), which is the number of cancer types in this study as shown in the following, then there are $M_0 L$ P -values (P_{ℓ}^m or \tilde{P}_{ℓ}^m), $P_h \in \mathbb{R}^{M_0 L}$ which are ranked in ascending order (i.e., if $h > h'$ then $P_h > P_{h'}$). P_h s re-ranked from P_{ℓ}^m and \tilde{P}_{ℓ}^m are compared with a paired Wilcoxon test with the alternative hypothesis that the P_h found

via SVD (i.e., re-ranked P_ℓ^m) is greater (i.e., less significant) than that given by HOSVD (i.e., re-ranked \tilde{P}_ℓ^m).

The coincidence with class labels can be evaluated as follows. $v_{\ell j}$ as well as $\tilde{v}_{\ell_1 j}$ can be regarded as the individual genes' representative profiles that represent dependence upon j . Thus, if $v_{\ell j}$ as well as $\tilde{v}_{\ell_1 j}$ have a dependence upon class labels, we can regard that individual genes' profiles have sufficient projection onto those coincident with class labels as well. Moreover, since $v_{\ell j}$ or $\tilde{v}_{\ell_1 j}$ can be computed by projecting x_{ij} onto $u'_{\ell i}$ or $\tilde{u}_{\ell_1 i}$, respectively, we can derive the dependence of x_{ij} upon j even only from the dependence upon i . In this sense, the evaluating coincident of $v_{\ell j}$ as well as $\tilde{v}_{\ell_1 j}$ with class labels can be regarded as a measure of inherent coincidence between individual genes' profiles and class labels as well.

2.3. Identification of Genes Expressed Distinctly between Class Labels and Enrichment Analysis

First, for each cancer dataset, using one of five class labels (see below), ℓ or ℓ_1 with the smallest P_ℓ^m or \tilde{P}_ℓ^m computed by categorical regression of Equations (6) and (7), respectively, we attributed p -values to individual genes (proteins) as follows:

$$P_i^m = P_{\chi^2} \left[> \left(\frac{u_{\ell i}}{\sigma_\ell} \right) \right] \tag{8}$$

or

$$\tilde{P}_i^m = P_{\chi^2} \left[> \left(\frac{\tilde{u}_{\ell_1 i}}{\sigma_{\ell_1}} \right) \right] \tag{9}$$

where $P_{\chi^2} [> x]$ is the cumulative χ^2 distribution when the argument is larger than x and σ_ℓ and σ_{ℓ_1} are optimized standard deviations so that P_i^m or \tilde{P}_i^m obeys Gaussian, which is the null hypothesis, as much as possible [10,11]. Computed p -values are corrected by the BH criterion [9] and i s associated with adjusted p -values less than 0.01 are selected. Selected genes in each of the 27 cancer types were separately uploaded to Enrichr [13] with the enrichR [14] package. Then, enrichment in KEGG, GO BP, GO CC, and GO MF were retrieved.

Enrichr evaluates the enrichment of genes using Fisher's exact test. Suppose G_1 is a set of genes uploaded (i.e., genes selected by our method) and G_2 is a set of genes with known function (e.g., genes that belong to a specific KEGG pathway). The overlap between G_1 and G_2 , $G_1 \cap G_2$, is evaluated by the comparison with that by chance. If $G_1 \cap G_2$ is much larger than that by chance and the probability of occurrence by chance is small, G_1 is regarded to be associated with the function associated with G_2 . In this study, we selected biological terms enriched if associated adjusted p -values given by Enrichr are less than 0.05.

For the enrichment analyses that considered HALLMARK cancer gene sets, we evaluate enrichment with enrichment function in clusterProfiler package [15] with the reference to HALLMARK genes in msigdb package [16]. Gene sets associated with adjusted p -values less than 0.05 were selected.

2.4. PPI Dataset

We have employed the following two PPI datasets for comparison.

2.4.1. Stanford PPI Dataset

The PPI dataset, PP-Pathways_ppi.csv.gz, was retrieved from the human protein-protein interaction network [17], which includes 342,353 pairs for 21,557 proteins. After excluding self-pairs (i.e., self-dimers), these data were formatted as a $n_{ij} \in \mathbb{R}^{N \times N}$ where $N = 21,557$. Then, there were 16,774 i s, which is common with the i s in the TCGA gene expression profiles (see below). Since these 342,353 pairs represent only one of n_{ij} and n_{ji} , (i.e., pairs whose order is reversed are not included), when $n_{ij} \neq n_{ji}$, $n_{ij} = n_{ji} = 1$ is required.

2.4.2. BIOGRID PPI Dataset

The PPI dataset, BIOGRID-MV-Physical-4.4.221.tab2.zip, which is supposed to represent physical PPIs, was retrieved from BIOGRID [18], which includes 437,679 pairs for 27,978 proteins. These data were also formatted as a $n_{ij} \in \mathbb{R}^{N \times N}$ where $N = 27,978$. Then, there were 11,294 *is*, which is common with the *is* in the TCGA gene expression profiles (see below). Since these 437,679 pairs represent only one of n_{ij} and n_{ji} , (i.e., pairs whose order is reversed are not included), when $n_{ij} \neq n_{ji}$, $n_{ij} = n_{ji}$ is required. In the BIOGRID dataset, n_{ij} is not always taken to be 1 when protein pairs interact with each other, but the number of occurrences in the BIOGRID PPI datasets. Thus, n_{ij} can be larger than 1 in contrast to the Stanford PPI. Thus, n_{ij} in the BIOGRID PPI represents not only whether pairs of proteins interact with each other, but also the strength of the interaction.

2.5. Gene Expression Profiles

TCGA gene expression profiles are used as x_{ij} . The RTCGA dataset [19] was used for this purpose. RTCGA.rnaseq [20] was used as a gene expression profile. It includes 27 cancer datasets with various sample sizes (*js*) ranging from a few tens to a few hundred, as well as 20,532 genes (*is*) whose expression profiles are available. The cancers considered are ACC, BLCA, BRCA, CESC, COAD, ESCA, GBM, HNSC, KICH, KIRC, KIRP, LGG, LIHC, LUAD, LUSC, OV, PAAD, PCPG, PRAD, READ, SARC, SKCM, STAD, TGCT, THCA, UCEC, and UCS. The class labels considered, retrieved from RTCGA.clinical [21], are patient.vital_status, patient.stage_event.pathologic_stage, patient.stage_event.tnm_categories, pathologic_categories, pathologic_m, patient.stage_event.tnm_categories.pathologic_categories, pathologic_n, and patient.stage_event.tnm_categories.pathologic_categories.pathologic_t. In order to avoid complexity, in the following, we employ shortened label class names as follows: “vital_status”, “pathologic_stage”, “pathologic_m”, “pathologic_n”, and “pathologic_t”. All 27 cancer datasets are associated with “vital_status” labels, “pathologic_stage” and “pathologic_m” are associated with only 18 datasets, and “pathologic_t” and “pathologic_n” are associated with 20 datasets (Table 1).

Table 1. Availability of class labels for cancer datasets. (1) “patient.vital_status”, (2) “pathologic_stage,” (3) “pathologic_m”, (4) “pathologic_t”, and (5) “pathologic_n”.

	ACC	BLCA	BRCA	CESC	COAD	ESCA	GBM	HNSC	KICH	KIRC	KIRP	LGG	LIHC	LUAD
(1)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
(2)	✓	✓	✓		✓	✓		✓	✓	✓	✓		✓	✓
(3)		✓	✓	✓	✓	✓		✓	✓	✓	✓		✓	✓
(4)	✓	✓	✓	✓	✓	✓		✓	✓	✓	✓		✓	✓
(5)	✓	✓	✓	✓	✓	✓		✓	✓	✓	✓		✓	✓
	LUSC	OV	PAAD	PCPG	PRAD	READ	SARC	SKCM	STAD	TGCT	THCA	UCEC	UCS	Total (= M_0)
(1)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	27
(2)	✓		✓			✓		✓	✓	✓	✓			18
(3)	✓		✓			✓		✓	✓	✓	✓			18
(4)	✓		✓		✓	✓		✓	✓	✓	✓			20
(5)	✓		✓		✓	✓		✓	✓	✓	✓			20

A total of 16,774 *is* (for Stanford) and 11,294 (for BIOGRID) that are common with *is* in PPI (see above) are used.

3. Results

3.1. Identification of Sample Vectors Coincident with Labels

The first evaluation of the proposed TB method is a comparison of the significance of coincidence between labels and vectors attributed to samples with and without consideration of PPI (i.e., comparisons between Equations (6) and (7)). If Equation (7) can provide

more significance than Equation (6), an integrated analysis of PPI and GE can improve the performance, because PPI itself is unlikely to include class label information, which is supposed to be specific to individual cancer types. The reason why we employed these class labels is simply because they are widely common for the majority of cancer types in TCGA. Since they are patient class labels, they might not be directly related to some specific biological concepts.

3.1.1. Stanford PPI

In this section, we evaluate each class label.

“vital_status”

First, we considered the label “vital_status”, which has two levels, “dead” and “alive.” Figure 2 (the left panel (1)) represents the logarithmic p -values computed by applying a Wilcoxon test (2.144×10^{-8}) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and $\tilde{v}_{\ell_{1j}}$ with Equation (7), respectively, whose scatter plot is shown in Figure 3 (the left panel). Since the number of vectors attributed to samples associated with adjusted p -values less than 0.05 for HOSVD (green and blue crosses) is larger than those for SVD (green crosses), the integrated analysis clearly improves the coincidence between the class label “vital_status” and vectors attributed to samples.

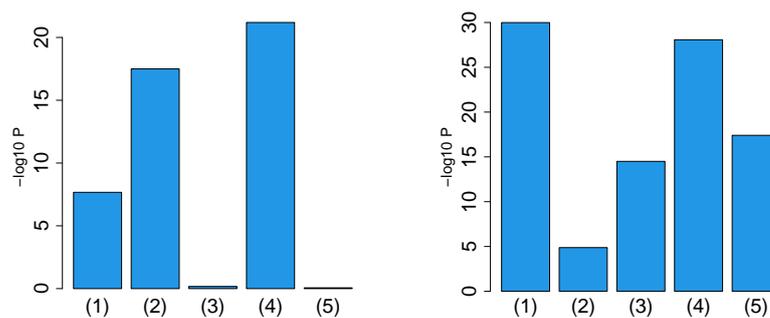


Figure 2. Barplot of p -values computed by a Wilcoxon test to evaluate the difference in ascending ordered P_h between SVD (Equation (6)) and HOSVD (Equation (7)) when Stanford PPI (left) or BIOGRID PPI (right) was used. (1) “vital_status”, (2) “pathologic_stage”, (3) “pathologic_m”, (4) “pathologic_t”, (5) “pathologic_n”. Numerical values of bar plots are listed in Table S1.

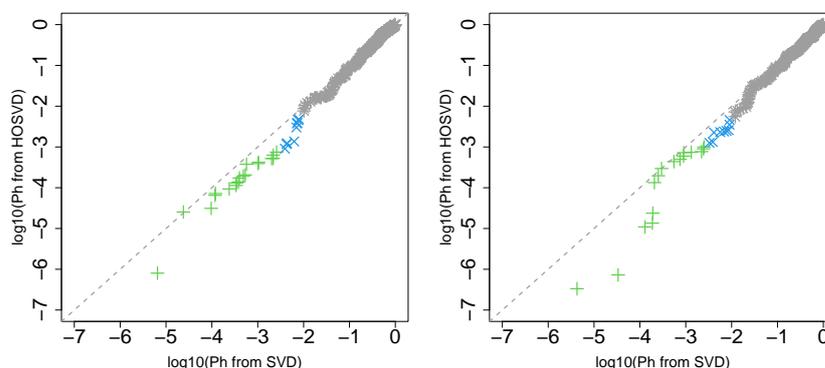


Figure 3. Scatter plot (logarithmic scale) of ascending ordered P_h computed from v_{ℓ_j} (horizontal axis) and $\tilde{v}_{\ell_{1j}}$ (vertical axis) for “vital_status.” Green crosses are those associated with adjusted p -values less than 0.05 for both axes and blue crosses are those associated with adjusted p -values less than 0.05 for the vertical axis alone. Grey asterisks represent all other situations. **Left:** Stanford PPI, **right:** BIOGRID PPI.

“pathologic_stage”

Next, we considered the label “pathologic_stage.” Figure 2 (the left panel (2)) represents the logarithmic p -values computed by applying a Wilcoxon test (3.178×10^{-18}) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and $\tilde{v}_{\ell_{1j}}$ with Equation (7), respectively, whose scatter plot is shown in Figure 4 (the left panel).

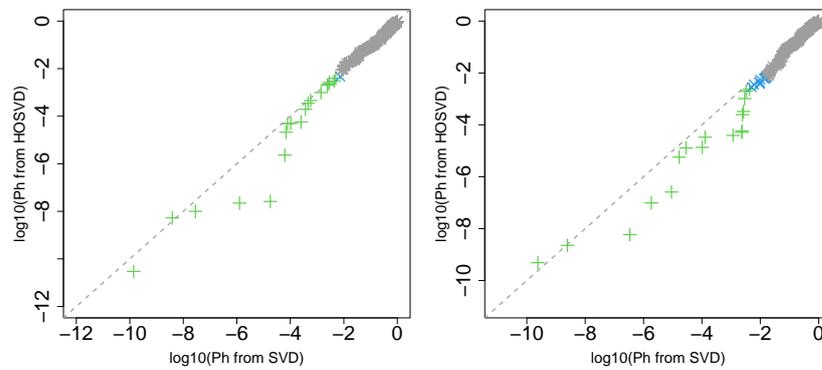


Figure 4. Scatter plot in logarithmic scale of ascending ordered P_h computed from v_{ℓ_j} (horizontal axis) and $\tilde{v}_{\ell_{1j}}$ (vertical axis) for “pathologic_stage.” Green crosses are those associated with adjusted p -values less than 0.05 for both axes and blue crosses are those associated with adjusted p -values less than 0.05 only for vertical axis. Grey asterisks represent all other situations. **Left:** Stanford PPI, **right:** BIOGRID PPI.

Again, P_h with integrated analysis of GE and PPI is significantly lower than that without consideration of PPI. Thus, the improvement observed in the label “vital_status” is unlikely accidental.

“pathologic_m”

Next, we considered the label “pathologic_m.” Figure 2 (the left panel (3)) represents the logarithmic p -values computed by applying a Wilcoxon test (0.6685) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and $\tilde{v}_{\ell_{1j}}$ with Equation (7), respectively, whose scatter plot is shown in Figure 5 (the left panel).

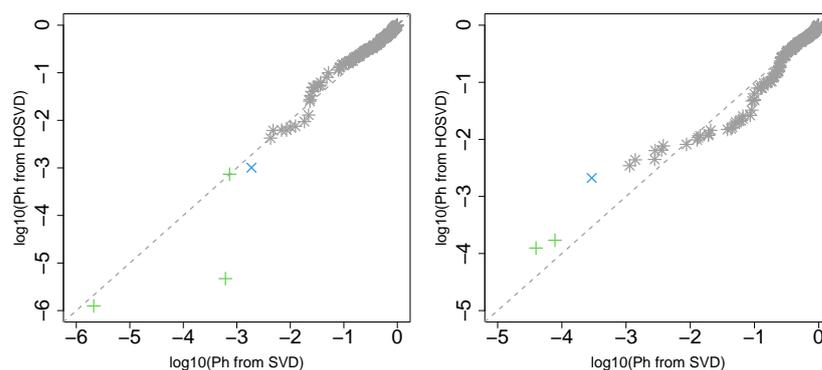


Figure 5. Scatter plot (logarithmic scale) of ascending ordered P_h computed from v_{ℓ_j} (horizontal axis) and $\tilde{v}_{\ell_{1j}}$ (vertical axis) for “pathologic_m.” Green crosses are associated with adjusted p -values less than 0.05 for both axes and blue crosses are associated with adjusted p -values less than 0.05 for the vertical axis alone. Grey asterisks represent all other situations. **Left:** Stanford PPI, **right:** BIOGRID PPI.

Since consideration of PPI does not improve the coincidence between class labels and vectors attributed to samples in this case, integrated analysis of PPI and GE does not always improve the coincidence (this will be discussed further).

“pathologic_t”

Next, we considered the label “pathologic_t.” Figure 2 (the left panel (4)) represents the logarithmic p -values computed by applying a Wilcoxon test (6.430×10^{-22}) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and $\tilde{v}_{\ell_{1j}}$ with Equation (7), respectively, whose scatter plot is shown in Figure 6 (the left panel).

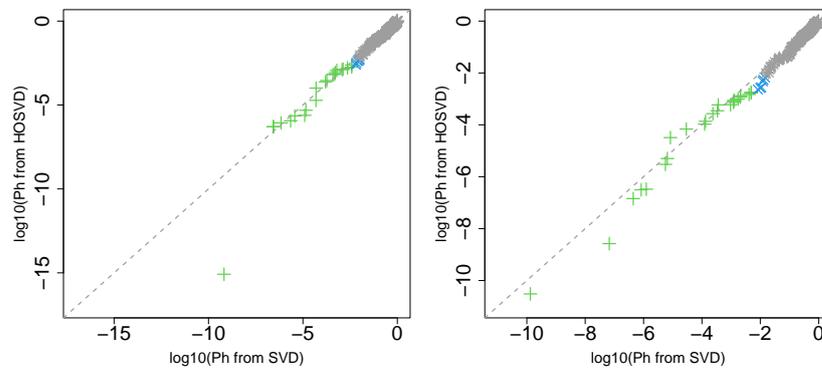


Figure 6. Scatter plot (logarithmic scale) of ascending ordered P_h computed from v_{ℓ_j} (horizontal axis) and $\tilde{v}_{\ell_{1j}}$ (vertical axis) for “pathologic_t.” Green crosses are associated with adjusted p -values less than 0.05 for both axes and blue crosses are associated with adjusted p -values less than 0.05 for the vertical axis alone. Grey asterisks represent all other situations. **Left:** Stanford PPI, **right:** BIOGRID PPI.

Although the $\log P_h$ does not appear significantly distinct between SVD and HOSVD (Figure 6, the left panel), since the p -values are small enough (Figure 2, the left panel (4)), integrated analysis of PPI and GE could improve the coincidence between vectors attributed to samples with class labels.

“pathologic_n”

Next, we considered the label “pathologic_n.” Figure 2 (the left panel (5)) represents the logarithmic p -values computed by applying a Wilcoxon test (0.8667) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and $\tilde{v}_{\ell_{1j}}$ with Equation (7), respectively, whose scatter plot is shown in Figure 7 (the left panel).

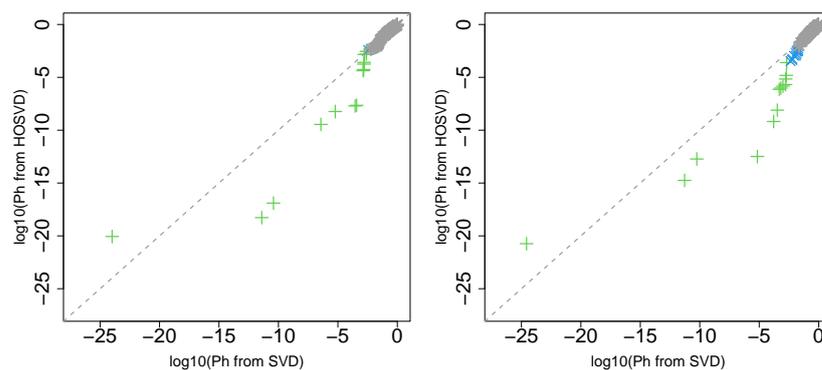


Figure 7. Scatter plot (logarithmic scale) of ascending ordered P_h computed from v_{ℓ_j} (horizontal axis) and $\tilde{v}_{\ell_{1j}}$ (vertical axis) for “pathologic_n.” Green crosses are associated with adjusted p -values less than 0.05 for both axes and blue crosses are associated with adjusted p -values less than 0.05 for the vertical axis alone. Grey asterisks represent all other situations. **Left:** Stanford PPI, **right:** BIOGRID PPI.

Since most vectors attributed to samples associated with adjusted P_h less than 0.05 for SVD are associated with lower P_h for HOSVD, integrated analysis of GE and PPI surely improved the coincidence between vectors attributed to samples and class labels.

3.1.2. BIOGRID PPI

Although integrated analysis of Stanford PPI and GE surely improved the coincidence between vectors attributed to samples and class labels for four out of five class labels (Figure 2), we were not certain if one class label, “pathologic_m”, without improved coincidence is because of PPI or because of the class label itself. Furthermore, we were not certain whether coincidence increased or decreased when we considered other PPIs. To examine these questions, we tested another PPI taken from BIOGRID. In the following section, we evaluate each class label.

“vital_status”

First, we considered the label “vital_status”, which has two levels, “dead” and “alive”. Figure 2 (the right panel (1)) represents the logarithmic p -values computed by applying a Wilcoxon test (1.027451×10^{-30}) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and \tilde{v}_{ℓ_j} with Equation (7), respectively, whose scatter plot is shown in Figure 3 (the right panel). Compared to the improvement when Stanford PPI is used (left panel in Figure 3), coincidence improvement increased when PPI is considered. This suggests that which PPI is used greatly affects the improvement when integrated analysis of GE and PPI is performed.

“pathologic_stage”

Next, we considered the label “pathologic_stage.” Figure 2 (the right panel (2)) represents the logarithmic p -values computed by applying a Wilcoxon test (1.34×10^{-5}) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and \tilde{v}_{ℓ_j} with Equation (7), respectively, whose scatter plot is shown in Figure 4 (the right panel). Compared to the left panel in Figure 4, coincidence improvement slightly decreased when PPI is considered. Although the direction of alteration is opposite to the “vital_status”, this suggests again that which PPI is used greatly affects the improvement when integrated analysis of GE and PPI is performed.

“pathologic_m”

Next, we considered the label “pathologic_m.” Figure 2 (the right panel (3)) represents the logarithmic p -values computed by applying a Wilcoxon test (3.16×10^{-15}) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and \tilde{v}_{ℓ_j} with Equation (7), respectively, whose scatter plot is shown in Figure 5 (the right panel). Although p -values computed by the Wilcoxon test are small enough to be significant, because no P_h s associated with significant adjusted p -values (the green and blue crosses in the right panels in Figure 5) decreased because of integrated analysis of PPI and GE, it is not regarded as an improvement. Thus, the failure of “pathologic_m” when Stanford PPI was used (left panel in Figure 5) is likely because of the class label itself and not because of PPI.

“pathologic_t”

Next, we considered the label “pathologic_t.” Figure 2 (the right panel (4)) represents the logarithmic p -values computed by applying a Wilcoxon test (8.534979×10^{-29}) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and \tilde{v}_{ℓ_j} with Equation (7), respectively, whose scatter plot is shown in Figure 6 (the right panel). Although two panels in Figure 6 do not look distinct, since p -values computed by the Wilcoxon test (Figure 2 (4) in the right panel) are much smaller than those when Stanford PPI was employed (Figure 2 (4) in the left panel), integrated analysis of PPI and GE could improve the coincidence between vectors attributed to samples with class labels.

“pathologic_n”

Next, we considered the label “pathologic_n.” Figure 2 (the right panel (5)) represents the logarithmic p -values computed by applying a Wilcoxon test (3.930971×10^{-18}) to ascending ordered $\log_{10} P_h$ computed from v_{ℓ_j} with Equation (6) and \tilde{v}_{ℓ_j} with Equation (7),

respectively, whose scatter plot is shown in the right panel of Figure 7 (the right panel). Because the right panel of Figure 7 is somewhat improved, compared to the left panel of Figure 7, the employment of BIOGRID PPI could improve the performance with Stanford PPI.

3.2. Identification of DEGs and Enrichment Analysis

In general, integrated analysis of PPI and GE could improve the coincidence between vectors attributed to samples and class labels. Nevertheless, it is still unclear whether the improved coincidence between vectors attributed to samples and class labels is useful. To address this problem, we performed an enrichment analysis of DEGs.

3.2.1. Stanford PPI

Figure 8 shows the summation of the number of biological terms enriched over 27 cancer classes for four categories (KEGG, GO BP, GO CC, and GO MF). A more detailed cancer cell barplot is available as supplementary material. It is obvious that integrated analysis of PPI and GE increases the number of enriched biological terms, no matter which category is considered.

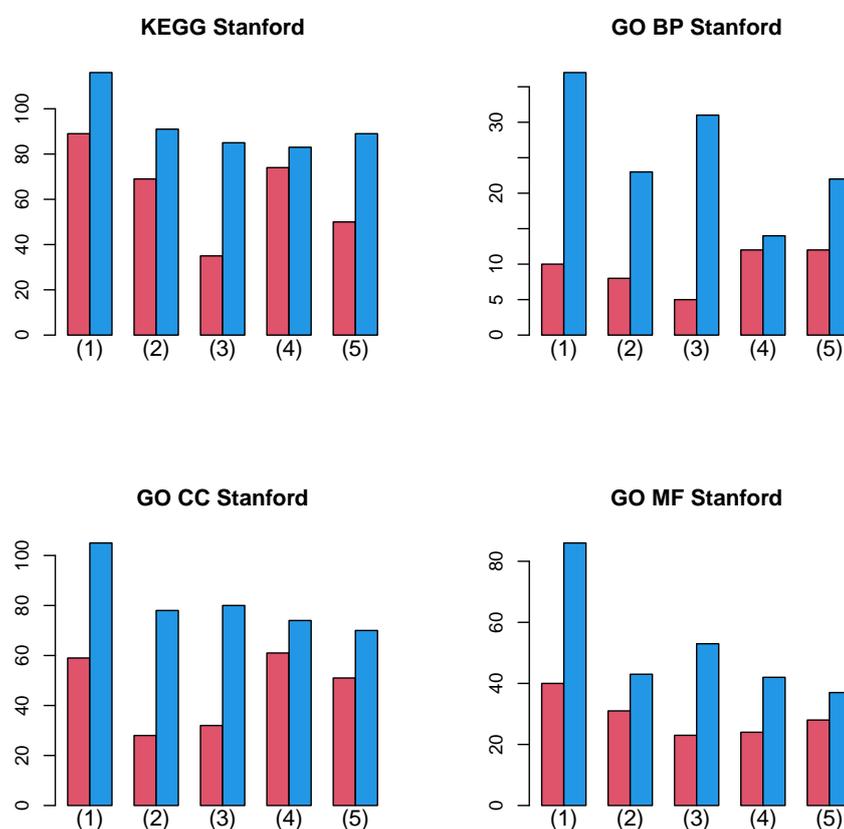


Figure 8. Barplot of the number of enriched biological terms summed over 27 cancers when Stanford PPI was used. (1) “vital_status”, (2) “pathologic_stage”, (3) “pathologic_m”, (4) “pathologic_t”, and (5) “pathologic_n”. Red: without integration of PPI, blue: with integration of PPI.

3.2.2. BIOGRID PPI

Figure 9 shows the summation of the number of biological terms enriched over 27 cancer classes for four categories (KEGG, GO BP, GO CC, and GO MF). A more detailed cancer cell barplot is available as supplementary material. It is obvious that the integrated analysis of PPI and GE increases the number of enriched biological terms excluding a few cases, although the amount of the increase is less than in Figure 8.

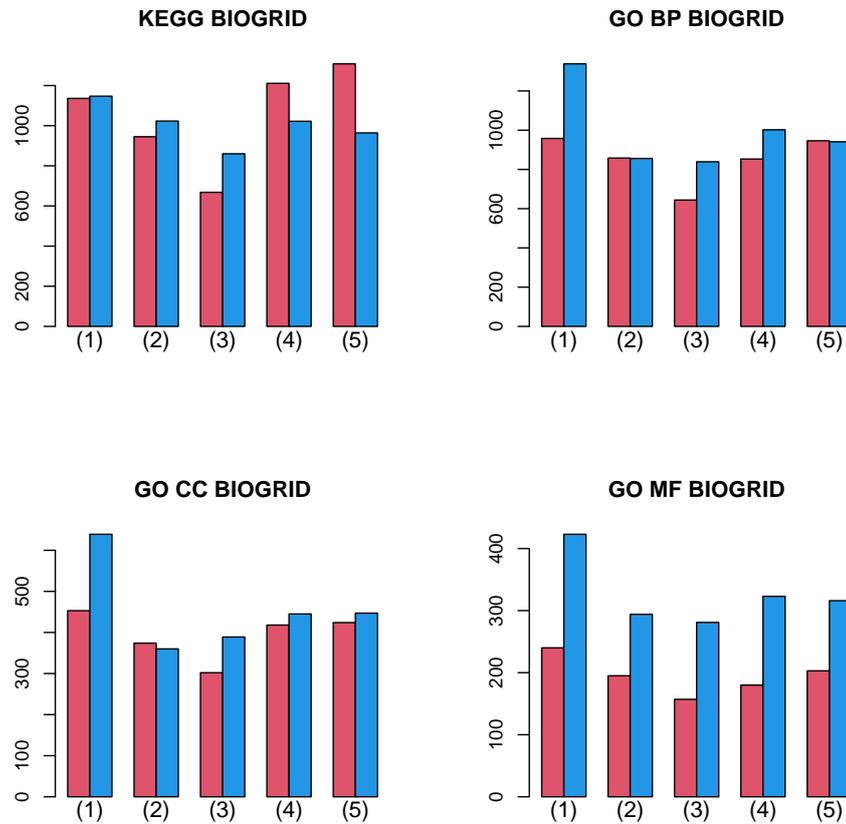


Figure 9. Barplot of the number of enriched biological terms summed over 27 cancers when BIOGRID PPI was used. (1) “vital_status”, (2) “pathologic_stage”, (3) “pathologic_m”, (4) “pathologic_t”, and (5) “pathologic_n”. Red: without integration of PPI, blue: with integration of PPI.

In conclusion, the integrated analysis of PPI and GE increases not only the coincidence between vectors attributed to samples and class labels but also the biological reliability of selected genes.

4. Discussion

Integrated analysis of PPI and GE is unlikely to improve coincidence with class labels because there are many class labels which are distinct from one another, whereas the PPI network does not vary depending on class labels. GE also does not change its values depending on class labels, but because samples are associated with class labels, it is not surprising that GE is coincident with class labels to some extent. Nevertheless, why can PPI that does not have any direct relation to class labels improve coincidence with class labels?

To clarify this point, we tried a simpler integration between PPI and GE. We computed sample vectors directly from PPI, not passing through TD. Mathematically, vectors attributed to samples can be computed directly from PPI as

$$\hat{v}_{\ell j} = \sum_{i=1}^N u_{\ell i} x_{ij} \tag{10}$$

where $u_{\ell i}$ was computed from PPI, n_{ii} , with Equation (1). Since x_{ij} is GE and $u_{\ell i}$ comes from PPI, it is a type of integrated analysis of GE and PPI. Then, coincidence with class labels is evaluated using categorical regression as

$$\hat{v}_{\ell j} = a''_{\ell} + \sum_s b''_{\ell s} \delta_{js} \tag{11}$$

which was just performed in the above analysis.

Next, we investigated the correlation between $\log_{10}\left(\frac{\tilde{P}_\ell^m}{P_\ell^m}\right)$ and $\log_{10}\left(\frac{\hat{P}_\ell^m}{P_\ell^m}\right)$ over 27 cancers, where \tilde{P}_ℓ^m , P_ℓ^m , and \hat{P}_ℓ^m are p -values computed from Equations (6), (7), and (11) for m th cancer type. This correlation evaluated whether the coincidence improvement between vectors attributed to samples and class labels by integrated analysis with TD, $\log_{10}\left(\frac{\tilde{P}_\ell^m}{P_\ell^m}\right)$, correlates with the coincidence improvement of the integrated analysis without TD evaluated by Equation (11) from coincidence in GE only evaluated by Equation (6), $\log_{10}\left(\frac{\hat{P}_\ell^m}{P_\ell^m}\right)$. Interestingly, regardless of class labels and PPI, there are significant correlations (Figure 10).

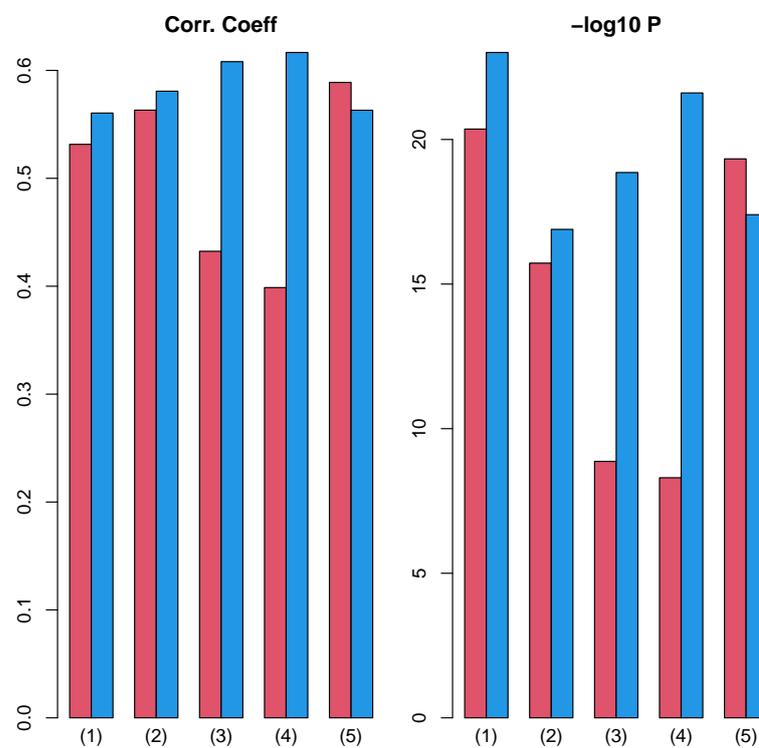


Figure 10. Left: Pearson correlation coefficient between $\log_{10}\left(\frac{\tilde{P}_\ell}{P_\ell}\right)$ and $\log_{10}\left(\frac{\hat{P}_\ell}{P_\ell}\right)$ over 27 cancers. Right: associated p -values (logarithmic scale) (1) “vital_status”, (2) “pathologic_stage”, (3) “pathologic_m”, (4) “pathologic_t”, and (5) “pathologic_n”. Red: Stanford PPI, blue: BIOGRID PPI. Numerical values of bar plots are listed in Table S2.

This means that improved coincidence between vectors attributed to samples and class labels caused by integrated analysis of PPI and GE with TD can be observed even in a simpler integration of PPI and GE (Equation (10)) to some extent, and this is why TB integrated analysis of PPI and GE can improve the coincidence between vectors attributed to samples and class labels.

Since one might wonder if a simpler integration of PPI and GE (Equation (10)) is powerful enough to improve the coincidence between vectors attributed to samples and class labels even without TD, we evaluated (Figure 11) the improvement of the coincidence using a simpler integration of PPI and GE, as in Figure 2. Since it is obvious that simpler integration of PPI and GE with Equation (10) is less likely to improve the coincidence

between vectors attributed to samples and class labels than TB integration, TB integration is required to improve the coincidence between vectors attributed to samples and class labels.

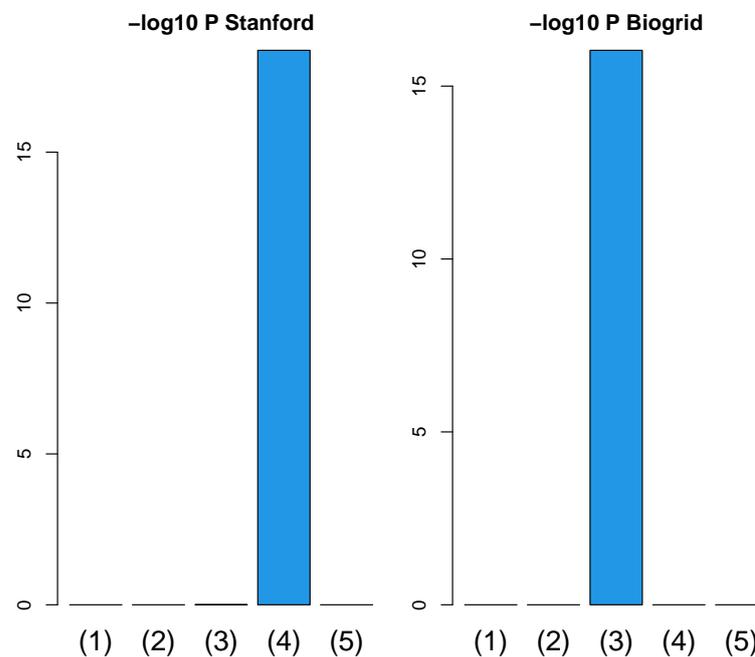


Figure 11. Barplot of p -values computed by a Wilcoxon test to evaluate the difference in ascending ordered P_h between SVD (Equation (6)) and simpler integration (Equation (11)). (1) "vital_status", (2) "pathologic_stage", (3) "pathologic_m", (4) "pathologic_t", and (5) "pathologic_n". **Left:** Stanford PPI, **right:** BIOGRID PPI. Numerical values of bar plots are listed in Table S3.

Considering the distinct performances of Stanford PPI (left panel in Figures 2 and 8), where n_{ij} takes only 1 or 0 dependent upon whether pairs of proteins interact or not and BIOGRID PPI (right panel in Figures 2 and 9), where n_{ij} can take larger values than 1 to represent the strength of interaction, it is likely better to consider not only if the interaction exists between pairs of proteins but also how strong the interaction between pairs of protein is. This finding might be able to help us to consider the integration of PPI and GE in the future.

We have also noted that the enrichment of SVD-based gene selection (i.e., without integration of PPI) in BIOGRID (the blue bars in Figure 9) is better than in Stanford PPI (the blue bars in Figure 8). This should not happen, since nothing can change between Stanford PPI and BIOGRID PPI when no integrated analyses were performed. Nevertheless, PPI can affect outcomes even before integration, since genes are screened based on whether they are also included in PPI. Because Stanford PPI and BIOGRID PPI differ, the other genes considered also differ. It turned out that this gene restriction largely affected the enrichment analysis. In this sense, integrated analysis can affect the outcome, but so does restricting the genes in considering overlaps with PPI.

To see if integrated analysis correctly identifies biologically reasonable genes, we investigated the frequency of selected terms in KEGG pathway for BIOGRID data (Table 2).

Table 2. Top ten most frequently selected pathways in KEGG pathway for BIOGRID data when PPI and GE are integrated. The numbers indicate the frequency of selections among 27 cancer types. (1) “vital_status”, (2) “pathologic_stage”, (3) “pathologic_m”, (4) “pathologic_t”, and (5) “pathologic_n”.

Pathway	(1)	(2)	(3)	(4)	(5)
Salmonella infection	25	18	17	20	20
JAK-STAT signaling pathway	23	18	17	19	20
Cytokine-cytokine receptor interaction	23	17	17	19	20
Influenza A	22	17	16	18	19
Pathways in cancer	22	17	14	19	19
Apoptosis	20	18	14	18	16
Ribosome	20	15	16	16	15
Non-alcoholic fatty liver disease	18	16	14	18	16
PI3K-Akt signaling pathway	16	17	15	18	16
C-type lectin receptor signaling pathway	18	15	13	17	16

Pathways in Table 2 are related to cancers. “Pathways in cancer” is directly related to cancer. “JAK-STAT signaling pathway” is related to cancers [22]. As for “Cytokine-cytokine receptor interaction”, Cytokine signaling is known to be related to cancers [23]. “PI3K-Akt signaling pathway” is known to be dysregulated almost in all human cancers [24]. As for “C-type lectin receptor signaling pathway”, C-Type lectin receptors are related to cancer immunity [25]. In addition to these, increased colon cancer risk was observed after severe Salmonella infection [26], cancers are associated with less apoptosis [27], non-alcoholic fatty liver disease increases the cancer risk [28], and influenza is associated with worse in-hospital clinical outcomes among hospitalized patients with malignancy [29]. Thus, most of the frequently selected pathways are related to cancers.

Table 2 included many biological pathways other than cancer-specific ones. To see more cancer-specific results, we restrict to “HALLMARK cancer gene sets” and repeated the above procedure for the integration with BIOGRID (Table 3).

Table 3. Enriched gene sets in “HALLMARK cancer gene sets” with and without the integration using BIOGRID. The numbers indicate the frequency of selections among 27 cancer types. (1) “vital_status”, (2) “pathologic_stage”, (3) “pathologic_m”, (4) “pathologic_t”, and (5) “pathologic_n”.

Pathway	(1)	(2)	(3)	(4)	(5)
HOSVD (with integration)					
HALLMARK_MYC_TARGETS_V1	10	4	5	6	6
HALLMARK_FATTY_ACID_METABOLISM	6	4	6	5	4
HALLMARK_OXIDATIVE_PHOSPHORYLATION	3	1	1	2	2
HALLMARK_APOPTOSIS	3	1	2	—	2
HALLMARK_PEROXISOME	2	2	—	1	—
HALLMARK_IL6_JAK_STAT3_SIGNALING	1	—	1	—	—
HALLMARK_MYC_TARGETS_V2	—	—	1	—	1
HALLMARK_ALLOGRAFT_REJECTION	1	—	—	—	—
HALLMARK_APICAL_SURFACE	1	—	—	—	—
SVD (without integration)					
HALLMARK_APOPTOSIS	11	5	7	3	3
HALLMARK_FATTY_ACID_METABOLISM	6	3	4	3	3
HALLMARK_PEROXISOME	2	2	—	3	3
HALLMARK_MYC_TARGETS_V1	1	2	2	1	1

Although the frequency of selection decreased from Table 2, possibly because of a more specific (strict) evaluation of the relationship to cancers, Table 3 still includes a substantial

frequency of selection. Moreover, the number of selected gene sets is more in the integrated analysis (denoted as HOSVD) than without integrated analysis (denoted as SVD). Therefore, even if we consider more cancer-specific features, the integrated analysis of PPI and GE has a substantial number of identifications and more numbers than that without integration

5. Conclusions

We proposed the integrated analysis of PPI and GE with TD that can result in more coincidence between vectors attributed to samples and one in five class labels in an evaluation of 27 cancer types using RNA-seq data retrieved from TCGA. Enrichment in genes selected as expressed distinctly among class labels are also improved. Furthermore, it was found that the consideration of the strength of PPI as well as the restriction of genes to intersect between PPI and GE can drastically improve the coincidence.

Supplementary Materials: The following supporting information can be downloaded at: <https://www.mdpi.com/article/10.3390/math11173655/s1>, Table S1: Numerical values of bar plots shown in Figure 2. Table S2: Numerical values of bar plots shown in Figure 10. Table S3: Numerical values of bar plots shown in Figure 11. More detailed enrichment analyses of those shown in Figures 8 and 9. The R source code used to perform this analysis.

Author Contributions: Y.-H.T. planned the research and performed the analyses. Y.-H.T. and T.T. evaluated the results, discussions, and outcomes and wrote and reviewed the manuscript. All authors have read and agreed to the published version of the manuscript

Funding: This work is supported in part by funds from the Chuo University (TOKUTEI KADAI KENKYU).

Data Availability Statement: All data analyzed in this paper are available in GEO.

Conflicts of Interest: The authors declare no conflict of interest. The funders had no role in the design of this study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the results.

References

- Jalili, M.; Gebhardt, T.; Wolkenhauer, O.; Salehzadeh-Yazdi, A. Unveiling network-based functional features through integration of gene expression into protein networks. *Biochim. Biophys. Acta (BBA) Mol. Basis Dis.* **2018**, *1864*, 2349–2359. [[CrossRef](#)] [[PubMed](#)]
- Elbashir, M.K.; Mohammed, M.; Mwambi, H.; Omolo, B. Identification of Hub Genes Associated with Breast Cancer Using Integrated Gene Expression Data with Protein-Protein Interaction Network. *Appl. Sci.* **2023**, *13*, 2403. [[CrossRef](#)]
- Karimizadeh, E.; Sharifi-Zarchi, A.; Nikaein, H.; Salehi, S.; Salamatian, B.; Elmi, N.; Gharibdoost, F.; Mahmoudi, M. Analysis of gene expression profiles and protein–protein interaction networks in multiple tissues of systemic sclerosis. *BMC Med. Genom.* **2019**, *12*, 199. [[CrossRef](#)] [[PubMed](#)]
- Tian, L.; Chen, T.; Lu, J.; Yan, J.; Zhang, Y.; Qin, P.; Ding, S.; Zhou, Y. Integrated Protein-Protein Interaction and Weighted Gene Co-expression Network Analysis Uncover Three Key Genes in Hepatoblastoma. *Front. Cell Dev. Biol.* **2021**, *9*, 631982. [[CrossRef](#)]
- Wu, C.; Zhu, J.; Zhang, X. Integrating gene expression and protein–protein interaction network to prioritize cancer-associated genes. *BMC Bioinform.* **2012**, *13*, 182. [[CrossRef](#)]
- Ewing, R.M.; Chu, P.; Elisma, F.; Li, H.; Taylor, P.; Climie, S.; McBroom-Cerajewski, L.; Robinson, M.D.; O'Connor, L.; Li, M.; et al. Large-scale mapping of human protein–protein interactions by mass spectrometry. *Mol. Syst. Biol.* **2007**, *3*, 89. [[CrossRef](#)]
- Su, L.; Liu, G.; Guo, Y.; Zhang, X.; Zhu, X.; Wang, J. Integration of Protein-Protein Interaction Networks and Gene Expression Profiles Helps Detect Pancreatic Adenocarcinoma Candidate Genes. *Front. Genet.* **2022**, *13*, 854661. [[CrossRef](#)]
- Zhong, J.; Tang, C.; Peng, W.; Xie, M.; Sun, Y.; Tang, Q.; Xiao, Q.; Yang, J. A novel essential protein identification method based on PPI networks and gene expression data. *BMC Bioinform.* **2021**, *22*, 248. [[CrossRef](#)]
- Taguchi, Y.H. *Unsupervised Feature Extraction Applied to Bioinformatics*; Springer International Publishing: Cham, Switzerland, 2020. [[CrossRef](#)]
- Taguchi, Y.H.; Turki, T. Adapted tensor decomposition and PCA based unsupervised feature extraction select more biologically reasonable differentially expressed genes than conventional methods. *Sci. Rep.* **2022**, *12*, 17438. [[CrossRef](#)]
- Taguchi, Y.H.; Turki, T. Application note: TDbasedUFE and TDbasedUFEadv: Bioconductor packages to perform tensor decomposition based unsupervised feature extraction. *Front. Artif. Intell.* **2023**, *6*, 1237542.
- Nakerekanti, M.; Narasimha, V. Analysis on Malware Issues in Online Social Networking Sites (SNS). In Proceedings of the 2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS), Coimbatore, India, 15–16 March 2019; pp. 335–338. [[CrossRef](#)]

13. Xie, Z.; Bailey, A.; Kuleshov, M.V.; Clarke, D.J.B.; Evangelista, J.E.; Jenkins, S.L.; Lachmann, A.; Wojciechowicz, M.L.; Kropiwnicki, E.; Jagodnik, K.M.; et al. Gene Set Knowledge Discovery with Enrichr. *Curr. Protoc.* **2021**, *1*, e90. [[CrossRef](#)]
14. Jawaid, W. *enrichR: Provides an R Interface to 'Enrichr'*, R Package Version 3.2; 2023. Available online: <https://cran.r-project.org/web/packages/enrichR/> (accessed on 23 August 2023)
15. Wu, T.; Hu, E.; Xu, S.; Chen, M.; Guo, P.; Dai, Z.; Feng, T.; Zhou, L.; Tang, W.; Zhan, L.; et al. clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. *Innovation* **2021**, *2*, 100141. [[CrossRef](#)]
16. Dolgalev, I. *msigdb: MSigDB Gene Sets for Multiple Organisms in a Tidy Data Format*, R Package Version 7.5.1; 2022. Available online: <https://cloud.r-project.org/web/packages/msigdb/msigdb.pdf> (accessed on 23 August 2023).
17. Human Protein-Protein Interaction Network. 2018. Available online: <https://snap.stanford.edu/biodata/datasets/10000/10000-PP-Pathways.html> (accessed on 23 August 2023).
18. Oughtred, R.; Rust, J.; Chang, C.; Breitkreutz, B.J.; Stark, C.; Willems, A.; Boucher, L.; Leung, G.; Kolas, N.; Zhang, F.; et al. The BioGRID database: A comprehensive biomedical resource of curated protein, genetic, and chemical interactions. *Protein Sci.* **2021**, *30*, 187–200. [[CrossRef](#)] [[PubMed](#)]
19. Kosinski, M. *RTCGA: The Cancer Genome Atlas Data Integration*, R Package Version 1.28.0; 2022. Available online: <https://rtcga.github.io/RTCGA/> (accessed on 23 August 2023).
20. Kosinski, M. *RTCGA.rnaseq: Rna-Seq Datasets from the Cancer Genome Atlas Project*, R Package Version 20151101.28.0; 2022. Available online: <https://bioconductor.org/packages/release/data/experiment/html/RTCGA.rnaseq.html> (accessed on 23 August 2023).
21. Kosinski, M. *RTCGA.clinical: Clinical Datasets from The Cancer Genome Atlas Project*, R Package Version 20151101.28.0; 2022. Available online: <https://bioconductor.org/packages/release/data/experiment/html/RTCGA.clinical.html> (accessed on 23 August 2023).
22. Brooks, A.J.; Putoczki, T. JAK-STAT Signalling Pathway in Cancer. *Cancers* **2020**, *12*, 1971. [[CrossRef](#)]
23. Lee, M.; Rhee, I. Cytokine Signaling in Tumor Progression. *Immune Netw.* **2017**, *17*, 214. [[CrossRef](#)]
24. Yang, J.; Nie, J.; Ma, X.; Wei, Y.; Peng, Y.; Wei, X. Targeting PI3K in cancer: Mechanisms and advances in clinical trials. *Mol. Cancer* **2019**, *18*, 26. [[CrossRef](#)] [[PubMed](#)]
25. Yan, H.; Kamiya, T.; Suabjakyong, P.; Tsuji, N.M. Targeting C-Type Lectin Receptors for Cancer Immunity. *Front. Immunol.* **2015**, *6*, 408. [[CrossRef](#)]
26. Mughini-Gras, L.; Schaapveld, M.; Kramers, J.; Mooij, S.; Neeffjes-Borst, E.A.; Pelt, W.v.; Neeffjes, J. Increased colon cancer risk after severe Salmonella infection. *PLoS ONE* **2018**, *13*, e0189721. [[CrossRef](#)] [[PubMed](#)]
27. Wong, R.S. Apoptosis in cancer: From pathogenesis to treatment. *J. Exp. Clin. Cancer Res.* **2011**, *30*, 87. [[CrossRef](#)]
28. Mantovani, A.; Petracca, G.; Beatrice, G.; Csermely, A.; Tilg, H.; Byrne, C.D.; Targher, G. Non-alcoholic fatty liver disease and increased risk of incident extrahepatic cancers: A meta-analysis of observational cohort studies. *Gut* **2022**, *71*, 778–788. [[CrossRef](#)] [[PubMed](#)]
29. Li, J.; Zhang, D.; Sun, Z.; Bai, C.; Zhao, L. Influenza in hospitalised patients with malignancy: A propensity score matching analysis. *ESMO Open* **2020**, *5*, e000968. [[CrossRef](#)] [[PubMed](#)]

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.